# Cross-Validation

## Problem

Split the dataset $A = [a_1, a_2, \ldots, a_N]$ of $N$ objects, each of which belongs to one of the $M$ classes ($\forall i: 1 \le a_i \le M$), into $K$ disjoint parts. Each object must fall into exactly one part. Splitting must satisfy several constraints.

## Sizes balance constraint

Let cnt($s$) is the number of objects belonging to the $s$ part, then
$$\forall s,t: |\text{cnt}(s) - \text{cnt}(t)| \le 1$$

## Classes distribution balance constraint

Let cnt($s,c$) is the number of objects with class $c$ belonging to the $s$ part, then
$$\forall s,t,c: |\text{cnt}(s,c) - \text{cnt}(t,c)| \le 1$$

## Randomness of splitting constraint

All possible splitting should be equiprobable. To be sure in that, for the same input run splitting several times and check that:
1. For each object probability of being in any part is equally distributed.
2. For each pair of objects probability of being together in the same part is equally distributed.

## Example

Let $N = 10$, $M = 4$, $K = 3$, and
$$A = [a_0 = 1, a_1 = 1, a_2 = 1, a_3 = 1, a_4 = 2, a_5 = 2, a_6 = 2, a_7 = 3, a_8 = 3, a_9 = 4]$$
Then:
- The splitting $S_1 = [a_0, a_1, a_2]$, $S_2 = [a_3, a_4, a_5]$, $S_3 = [a_6, a_7, a_8, a_9]$ satisfy the first constraint but it isn't satisfy the second.
- The splitting $S_1 = [a_0, a_1, a_4, a_7, a_9]$, $S_2 = [a_2, a_5, a_8]$, $S_3 = [a_3, a_6]$ satisfy the second constraint but it isn't satisfy the first.
- The splitting $S_1 = [a_0, a_1, a_4, a_7]$, $S_2 = [a_2, a_5, a_8]$, $S_3 = [a_3, a_6, a_9]$ satisfy the both constraints.